

Integrating SAM Supervision for 3D Weakly Supervised Point Cloud Segmentation

Lechun You, Zhonghua Wu, Weide Liu, Xulei Yang, Jun Cheng, *Senior Member, IEEE*, Wei Zhou, *Senior Member, IEEE*, Bharadwaj Veeravalli, *Senior Member, IEEE*, Guosheng Lin

Abstract—Current methods for 3D semantic segmentation propose training models with limited annotations to address the difficulty of annotating large, irregular, and unordered 3D point cloud data. They usually focus on the 3D domain only, without leveraging the complementary nature of 2D and 3D data. Besides, some methods extend original labels or generate pseudo labels to guide the training, but they often fail to fully use these labels or address the noise within them. Meanwhile, the emergence of comprehensive and adaptable foundation models has offered effective solutions for segmenting 2D data. Leveraging this advancement, we present a novel approach that maximizes the utility of sparsely available 3D annotations by incorporating segmentation masks generated by 2D foundation models. We further propagate the 2D segmentation masks into the 3D space by establishing geometric correspondences between 3D scenes and 2D views. We extend the highly sparse annotations to encompass the areas delineated by 3D masks, thereby substantially augmenting the pool of available labels. Furthermore, we apply confidence- and uncertainty-based consistency regularization on augmentations of the 3D point cloud and select the reliable pseudo labels, which are further spread on the 3D masks to generate more labels. This innovative strategy bridges the gap between limited 3D annotations and the powerful capabilities of 2D foundation models, ultimately improving the performance of 3D weakly supervised segmentation.

Index Terms—Weakly Supervised Semantic Segmentation, 3D Point Cloud Segmentation, Scene Understanding

I. INTRODUCTION

3D semantic segmentation is a critical computer vision task that assigns semantic labels to each voxel in a three-dimensional scene or each point in a 3D point cloud. This facilitates a thorough comprehension and distinction

of various objects or structures in the scene. 3D semantic segmentation has a variety of applications such as autonomous driving, robotics, augmented reality, etc. These techniques can be categorized into three main classes: point-based methods [1]–[7], voxel-based methods [8], [9], and projection-based methods [10].

3D point cloud data is a collection of data points obtained by sensors such as LiDAR, RGB-D cameras, laser scanners, etc., with sparse, irregularly sampled, high-dimensional, ambiguous, and unordered characteristics. Due to the intrinsic characteristics of 3D data, annotating extensive data and acquiring high-quality labels poses many challenges. Therefore, many current approaches explore the use of sparse 3D annotations during training [11]–[18].

For weakly supervised semantic segmentation tasks, existing research often focuses on 3D point clouds independently, with limited efforts directed toward the integration with 2D images. However, 2D images, known for their simplicity, ease of capture, and detailed textures, naturally enhance the abundance of geometric data offered by 3D data. On the other hand, the convergence of 2D and 3D joint learning primarily revolves around fully supervised semantic segmentation tasks, where features of both modalities are projected and fused [19]–[21], or labels are propagated across modalities [22], with limited exploration in the domain of weakly supervised tasks.

Additionally, the advances in robust 2D foundational models with zero-shot capabilities have significantly matured the performance of 2D semantic segmentation. Leveraging this progress, we seek to exploit the geometric correspondence between 3D and 2D data, thereby maximizing the utility of sparsely available 3D annotations while still achieving strong performance compared to fully supervised methods. We introduce a novel weakly supervised semantic segmentation model designed to address the challenges posed by sparsely annotated 3D point cloud data. This model harnesses the strengths of 2D foundational models in image segmentation. When faced with unlabeled 2D images, we begin by employing the cutting-edge 2D semantic segmentation model, Semantic-SAM [23], to produce segmentation masks. We utilize the spatial regions without assigning specific semantic labels, thereby enabling adaptable alignment of class information. Then we leverage geometry information to back-project the 2D masks into the 3D domain and fuse the masks of the same object from different views, yielding 3D masks.

We employ diverse strategies to maximize the utility of

This research is partly supported by the Agency for Science, Technology and Research (A*STAR) under its MTC Programmatic Funds (Grant No. M23L7b0021).

L. You is with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore 119077 (email: lechun.you@u.nus.edu).

Z. Wu is with SenseTime Research, Singapore 069547 (email: wuzhonghua@sensetime.com).

W. Liu is with the School of Computer Science and Engineering, Nanyang Technological University, Singapore 639798 (email: weide001@e.ntu.edu.sg).

X. Yang and J. Cheng are with the Institute for Infocomm Research (I2R), Agency for Science, Technology and Research (A*STAR), Singapore 138632 (email: yang_xulei@i2r.a-star.edu.sg, cheng_jun@i2r.a-star.edu.sg).

W. Zhou is with School of Computer Science and Informatics, Cardiff University, UK (zhouw26@cardiff.ac.uk).

B. Veeravalli is with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore 119077 (email: elebv@nus.edu.sg).

G. Lin is with the School of Computer Science and Engineering, Nanyang Technological University, Singapore 639798 (email: gslin@ntu.edu.sg).

Manuscript received April 19, 2021; revised August 16, 2021.

limited annotations and 3D masks. Initially, we extend the original sparse annotations onto the 3D masks, significantly increasing the number of available labels. During training, we enforce consistency regularization [18] to categorize the predictions into reliable and ambiguous subsets and further distribute the reliable pseudo labels across specific regions delineated by the 3D masks in a proportional manner. These expanded labels are essentially accurate, thereby enhancing subsequent training iterations. Acknowledging their potential inadequacy in fully representing each masked region, as well as the presence of noise in the projected masks, we consider them noisy and apply a noise-robust loss [24] on them to mitigate this issue.

Our contributions can be summarized as follows.

- We present a novel approach for weakly supervised 3D semantic segmentation by utilizing 2D foundation models. By leveraging geometric correspondences between 2D and 3D data, we bridge the gap between 2D foundation models and 3D learning, significantly enhancing label availability by back-projecting 2D segmentation masks into 3D space and extending sparse annotations.
- We propose a strategy that fuses segmentation masks projected into 3D from different views. This is achieved by evaluating the overlap between different projected masks to determine whether they should be merged into a single mask, ultimately ensuring that each 3D object's mask is complete and non-redundant.
- We identify reliable pseudo labels and propose a strategy to refine and expand them using back-projected masks. This strategy is based on the proportion of pseudo labels within each projected mask, optimizing the threshold to achieve maximum accuracy. Our approach significantly increases the number of high-quality labels and enhances overall performance.
- Our approach achieves state-of-the-art performance, showcasing its effectiveness in harnessing minimal annotations for 3D data and unlabeled 2D data toward comprehensive 3D scene understanding.

II. RELATED WORK

A. 3D Semantic Segmentation

The research on 3D semantic segmentation is currently widespread, with a focus on leveraging the rich geometric information inherent in 3D data.

Charles et al. proposed PointNet [1], a pioneering work that directly processes irregular raw point cloud data, preserving permutation invariance by applying the symmetric function to aggregate global information. The latter PointNet++ [2] addresses its issue of insufficient capability to capture local information by utilizing a hierarchical structure. Moreover, Zhao et al. introduced the Point Transformer [6], which utilizes self-attention networks in 3D point cloud processing. Wu et al. further improved the Point Transformer by introducing group vector attention and position encoding multiplier [7]. Furthermore, the original sample-based pooling is simplified to partition-based pooling. To enhance efficiency and accuracy,

they further introduced Point Transformer V3 [25]. This version replaces the traditional neighbor search from K-Nearest Neighbors (KNN) with serialized neighbor mapping, streamlines attention patch interaction mechanisms, and simplifies positional encoding.

Different from handling points directly, Choy et al. [9] represented the point cloud data as sparse tensors and introduced MinkowskiNet with generalized sparse convolutions for processing 3D sparse voxel grids. Another kind of approach projects the 3D point cloud to multiple 2D views and takes advantage of the capability of 2D CNN and the fine-grained texture of 2D images [10].

Some current methods combine 2D and 3D data in integrated learning schemes. Kundu et al. rendered synthetic 2D images from virtual views of the 3D scene to train a 2D segmentation model, then generated predictive features for fusion on 3D mesh vertices [19]. To enable joint reasoning for 2D and 3D, Hu et al. proposed the Bidirectional Projection Network (BPNNet) [21], consisting of symmetric 2D and 3D subnets connected at the same decoder level. With a link matrix, 2D and 3D features are correspondingly projected to each other and fused.

However, these methods require fully annotated data to realize their full potential. Yet, annotating point cloud data is resource-intensive due to its unordered and irregular nature.

B. Weakly Supervised 3D Semantic Segmentation

To address the time-consuming and labor-intensive issue of annotating 3D point cloud data, some methods have proposed training semantic segmentation models on a limited set of labels. Liu et al. introduced the ‘One-Thing-One-Click (OTOC)’ scheme [13], which means that each object needs only one annotated label. It conducts iterative training and pseudo label propagation, supported by a relation network that learns the similarity among voxels. Nonetheless, the generated pseudo labels are noisy, and the ‘OTOC’ annotation remains costly, as the user must identify every individual object in the scene, with each scene requiring up to 2 minutes of annotation.

Some current research utilizes consistency regularization to weakly supervised point cloud segmentation. Wu et al. introduced the PointMatch [17] approach that utilizes point-wise predictions from one view as pseudo labels for another view within a point cloud scene, enhancing the overall consistency between the views. Wu et al. proposed RAC-Net [18], which leverages both prediction confidence and model uncertainty to divide the pseudo labels into the reliable set and the unreliable set. Then, different consistency constraints are applied to them. These methods directly apply Cross-Entropy loss to pseudo labels and predicted values, but they do not explore further utilization of these pseudo labels.

ActiveST [26] employs an active learning approach for weakly supervised point cloud segmentation tasks. This method automatically selects points with high potential for improving the model, based on prediction uncertainty, for manual annotation, while also using highly confident predictions as pseudo labels for training. However, it requires the user to manually annotate the selected points in each iteration,

which distinguishes it from other traditional weakly supervised training methods.

Dong et al. [27] utilized the Segment Anything Model (SAM) [28] to generate segmentation masks for 2D views and propagated sparse 2D labels on these masks. They projected the expanded labels from each view into the 3D scene and employed a voting strategy to aggregate these labels. However, they only utilized the masks in the 2D plane without further exploiting them in the 3D space, and thus missed opportunities to fully leverage the masks' potential in 3D applications. Instead of projecting the expanded 2D labels, we project the 2D segmentation masks to 3D masks to enable flexible use, so that more 3D labels can be expanded on them.

C. 2D Foundation Models

2D semantic segmentation [29]–[32] is to label each pixel in an image with a semantic class to provide a detailed understanding of the content of the image. Long et al. proposed Fully Convolutional Networks (FCN) [33] for segmentation and introduced a skip architecture that integrates semantic details from a deep, coarse layer with appearance features from a shallow, fine layer, which has been widely adopted as the backbone in many previous works [34], [35]. In 2021, Dosovitskiy et al. proposed the Vision Transformer (ViT) [36], which utilizes self-attention mechanisms. It divides an image into fixed-size patches and preserves spatial information through positional embeddings. ViT exhibits strong scalability and high performance, which makes it widely adopted in various image-processing tasks. Currently, many methods also improve upon ViT for 2D image segmentation [37], [38].

In recent times, with the rapid advancement of Large Language Models (LLM), some studies have integrated open vocabulary into visual tasks to establish foundational models for image segmentation. They seek to divide an image into semantic regions based on arbitrary text descriptions. Kirillov et al. proposed the Segment Anything Model (SAM) [28], which supports any segmentation prompt with zero-shot generalization. SAM3D [39] utilized SAM to generate masks for RGB images and then projected these masks onto 3D point clouds. By iteratively merging the masks, they achieved fine-grained segmentation results. However, this method strictly requires 2D views that are matched with the 3D scene. When inferring new data, images need to be re-collected for the new data. However, our approach only utilizes 2D images during the training phase to alleviate the annotation burden, while during inference, it directly uses 3D point clouds.

Building upon SAM, Li et al. introduced Semantic-SAM [23], which further enables image segmentation at any given granularity by jointly training the model on seven datasets with multiple semantic labels. Wang et al. introduced a hierarchical representation approach encompassing semantic, instance, and part levels [40]. This approach decouples the representation learning modules and text-image fusion mechanisms for both background and foreground.

2D foundation models, trained on large-scale 2D datasets and incorporating language modalities, demonstrate excellent generalization capabilities and zero-shot performance. This

makes them ideal complements to 3D data with the inherent geometric correspondences. We opt for the Semantic-SAM [23] to generate 2D semantic segmentation masks due to its strong performance and flexible granularity options, including semantic, instance, and part levels.

D. Learning From Noisy Labels

Several techniques have been proposed to address the challenge of training accurate deep-learning models in the presence of noisy labels. An approach to mitigate the impact of noisy labels is to use robust loss functions, which offer simplicity in implementation. Ghosh et al. theoretically demonstrated that loss functions such as the Mean Absolute Error (MAE) exhibit noise-robust properties, whereas the commonly employed Cross-Entropy (CE) loss function does not possess such robustness [41]. Zhang et al. introduced the Generalized Cross-Entropy (GCE) loss [42], which can be readily applied to existing networks. Furthermore, Wang et al. augmented the original CE loss with a Reverse Cross-Entropy (RCE) term, forming the Symmetric Cross-Entropy (SCE) loss [43]. Furthermore, Ma et al. demonstrated that normalization could render any loss robust to noisy labels [24]. They compared the performance of models trained with different combinations of normalized loss functions. Zhou et al. rectified commonly used loss functions and proposed asymmetric loss functions to deal with multiple types of noise [44].

In addition to utilizing robust loss functions, other methodologies focus on detecting and rectifying mislabeled data points [45]. For the point cloud segmentation task, Ye et al. proposed a confidence-based approach to select reliable labels, using historical predictions for each data point, and employing a voting strategy to generate labels within clusters [46].

III. METHOD

Our weakly supervised setting is defined as follows: during training, the required data includes 3D point clouds X , where only a small subset of points are annotated, referred to as limited annotations Y , while the rest remain unlabeled, along with RGB images that correspond to each point cloud. During inference, only the 3D point cloud is required.

The entire framework is illustrated in Figure 1. First, Semantic-SAM [23] is used to generate 2D segmentation masks for each view within a scene. These masks accurately outline pixel clusters corresponding to each object. Subsequently, leveraging the spatial correspondence between 2D views and the 3D scene, these 2D masks are back-projected onto the 3D space, discussed in Section III-A. We explore various strategies to maximize the utility of the 3D masks. First, we extend the limited 3D annotations onto the masks, creating additional training labels. But they are still around the initial annotations, failing to cover all the masked regions. To make full use of the masks, we select reliable pseudo labels and propagate them onto 3D masks to expand their coverage during training. They are further fused with the initial expanded annotations, discussed in Section III-C. However, despite the increase in the number of labels, noise is introduced. To mitigate this issue, a noise-robust normalized

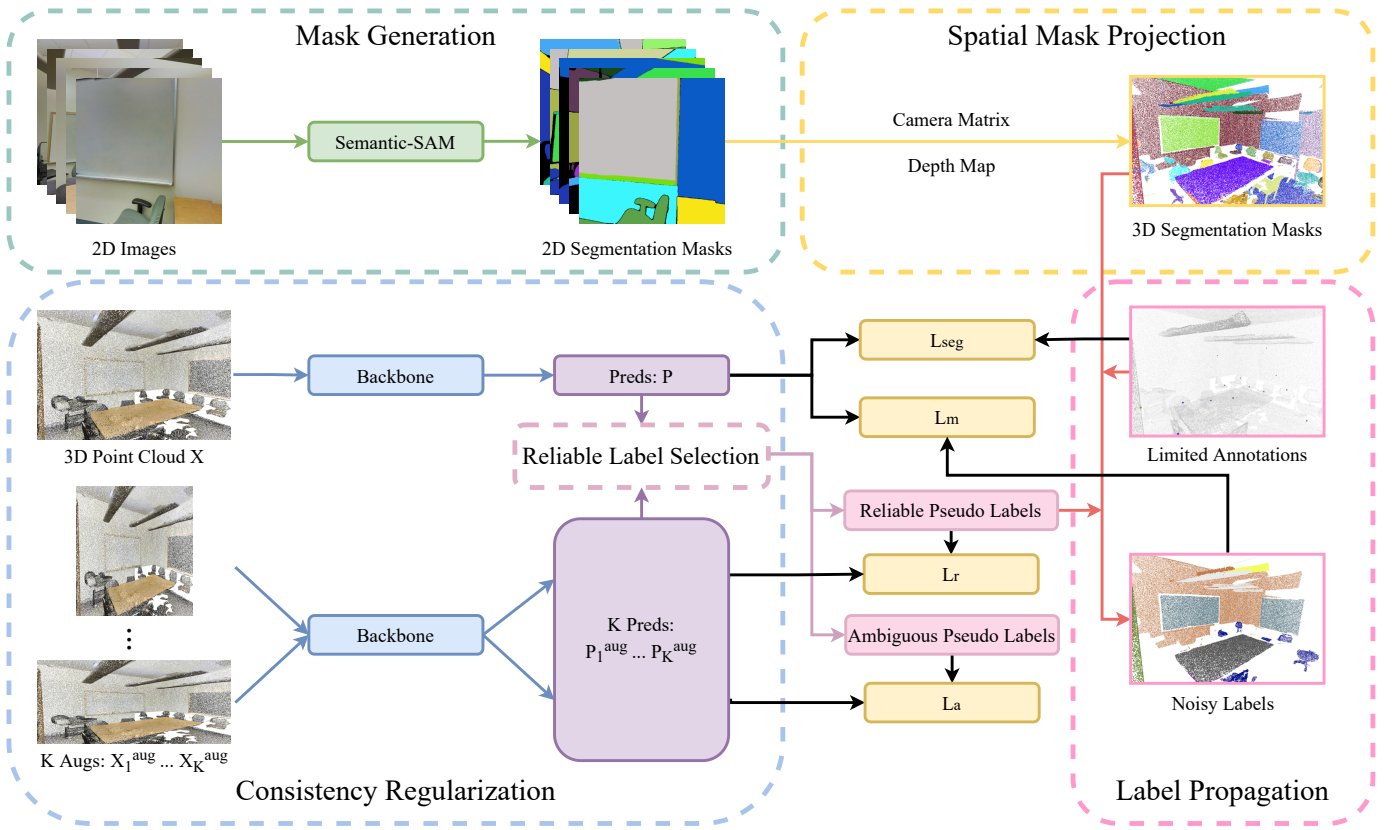


Fig. 1. Our method consists of four main components. The Mask Generation and Spatial Mask Projection modules are to generate segmentation masks for 2D views and project them into 3D space, as discussed in Section III-A. The Consistency Regularization module outlines the process of obtaining reliable pseudo labels and is detailed in Section III-B. The Label Propagation module is responsible for spreading limited annotations and reliable pseudo labels to masked regions, as elaborated in Section III-C.

loss [24] is implemented to train the model with the noisy labels, detailed in Section III-D. To classify the pseudo labels into reliable and ambiguous subsets, consistency regularization is employed based on prediction confidence and uncertainty [18], detailed in Section III-B.

A. Spatial Mask Projection

To harness the robust capabilities of foundation models for optimizing the utilization of limited annotations in our segmentation tasks, the Semantic-SAM [23] is employed to generate segmentation masks for 2D views. We chose Semantic-SAM because it has zero-shot capabilities, supports flexible segmentation at six different granularities, and allows for combining different granularities to generate segmentation masks.

For all 2D views corresponding to each 3D scene, adjacent views often have significant overlap, with the same object being captured in multiple frames. To save runtime and computational resources, we uniformly sample N_{view} views from all available 2D views to ensure that the captured objects are as complete as possible.

For each sampled 2D view, the set of segmentation masks is represented as a boolean matrix $Mask_{2D}$ of dimensions $[M, H, W]$, where M denotes the number of classes, H and W represent the height and width of the image. There are M boolean values corresponding to each pixel, with only one

being true, indicating that the pixel is assigned to the class represented by that index.

After that, $Mask_{2D}$ are propagated to the 3D domain using the spatial correspondence between 2D pixels and 3D points through a link matrix \mathcal{L} [21]. The initial step is to compute matrix \mathcal{M} , the multiplication of the intrinsic camera calibration matrix \mathbf{K} by the extrinsic camera pose matrix $[\mathbf{R}|\mathbf{t}]$ consisting of rotation \mathbf{R} and translation \mathbf{t} , denoted as:

$$\mathcal{M} = \mathbf{K}[\mathbf{R}|\mathbf{t}] \quad (1)$$

Therefore, the projection from 3D homogeneous coordinates $[x_i, y_i, z_i, 1]^T$ to 2D homogeneous coordinates $[u_i, v_i, 1]^T$ corresponding to the i^{th} 3D point can be denoted as:

$$[u_i, v_i, 1]^T = \mathcal{M}[x_i, y_i, z_i, 1]^T \quad (2)$$

The treatment of occlusion is addressed by incorporating depth information. Finally, \mathcal{L} is an $N \times 3$ matrix which reveals correspondences between 2D and 3D points:

$$\mathcal{L}_i = [u_i, v_i, m_i], \quad (3)$$

$$m_i = \begin{cases} 1, & \text{if } U_{min} \leq u_i \leq U_{max} \\ & \text{and } V_{min} \leq v_i \leq V_{max} \\ & \text{and } |d(u_i, v_i) - z'_i| \leq \delta \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

where N is the number of 3D points; m_i is the binary mask reveals whether the i^{th} 3D point has the valid projected 2D pixel; $U_{min}, U_{max}, V_{min}, V_{max}$ are the boundaries of the view frustum; $d(\cdot)$ is the mapping from coordinates to the depth; z'_i is the projected z coordinate of the point; δ is the threshold for depth matching.

For each 2D view, its segmentation masks $Mask_{2D}$ are projected onto the 3D masks $Mask'_{3D}$ based on the above correspondence. Therefore, the projected matrix $Mask'_{3D}$ has the dimensions of $[M, N]$, calculated by:

$$Mask'_{3D}(i) = \begin{cases} Mask_{2D}(u_i, v_i)_{M \times 1}, & \text{if } m_i = 1 \\ \text{False}_{M \times 1}, & \text{if } m_i = 0 \end{cases} \quad (5)$$

As a result, we obtained N_{view} sets of 3D masks, each corresponding to the masks from different 2D views. These 3D masks are then merged to handle overlapping regions. The merging process works as follows: we initialize $Mask_{3D}$ with the first set of 3D masks projected from the first view, and then sequentially examine the subsequent sets of 3D masks. If a mask from the current set (current view) overlaps with a mask from the previous set (previous view) by more than a certain threshold, we assume they belong to the same class and merge the two masks. Otherwise, the current mask is added as a new mask. This process continues until all masks have been processed. The final result is a matrix $Mask_{3D}$ with dimensions $[T, N]$, where T is the number of object classes after merging. Each 3D point is assigned T boolean values, indicating whether it belongs to one of the T classes.

B. Consistency Regularization and Reliable Pseudo Label Selection

Motivated by the methodology of RAC-Net [18], we seek to enhance the utilization of limited annotations and select reliable pseudo labels through the incorporation of consistency regularization.

We feed the original point cloud X into the backbone network (which can be any point cloud segmentation network), alongside its K augmented counterparts $X_{1...K}^{aug}$, to obtain P and $P_{1...K}^{aug}$. The augmentation methods we use are PointWOLF [47] and Affine Transformation (AT).

$$\begin{aligned} P &= \text{backbone}(X) \\ P_1^{aug} &= \text{backbone}(X_1^{aug}) \\ &\dots \\ P_K^{aug} &= \text{backbone}(X_K^{aug}) \end{aligned} \quad (6)$$

Subsequently, following the procedure employed in the RAC-Net, the predictions are categorized into a reliable set P^r and an ambiguous set P^a :

$$\begin{aligned} P^r &= R \cdot P, \quad P^a = (1 - R) \cdot P, \\ R &= \mathbb{1} \sum_{c=1}^C (\mathbb{1}[\bar{P}_c \geq \tau] \cdot \mathbb{1}[\sigma(\hat{P}_c) \leq \kappa]) > 0 \end{aligned} \quad (7)$$

where \bar{P}_c denotes the confidence of predictions, calculated from the mean of predictions for both the original and augmented data; $\sigma(\hat{P}_c)$ signifies the uncertainty of the predictions,

which is the statistical variance; $\mathbb{1}$ is the indicator function. For the prediction of each point, if both conditions are satisfied — the prediction confidence \bar{P}_c is greater than or equal to the threshold τ and the uncertainty $\sigma(\hat{P}_c)$ is less than or equal to the threshold κ — then the prediction is considered reliable; otherwise, it is considered unreliable.

The reliable predictions P^r are transformed into one-hot pseudo labels \tilde{Y}^r . Then, the Cross-Entropy loss is applied between \tilde{Y}^r and predictions of augmented data:

$$\begin{aligned} \mathbf{L}_r &= CE[\tilde{Y}^r, R \cdot P_1^{aug}] \\ &+ \dots + CE[\tilde{Y}^r, R \cdot P_K^{aug}] \end{aligned} \quad (8)$$

For the ambiguous predictions P^a , the KL Divergence is computed between the soft pseudo labels P^a and the predictions of augmented data:

$$\begin{aligned} \mathbf{L}_a &= KL[P^a, (1 - R) \cdot P_1^{aug}] \\ &+ \dots + KL[P^a, (1 - R) \cdot P_K^{aug}] \end{aligned} \quad (9)$$

C. Label Propagation

In this module, we propagate both the limited annotations Y and reliable pseudo labels \tilde{Y}^r onto the masks, which significantly increases the number of available labels.

Label Initialization. Given the constraint of very limited labels Y , more labels are produced by propagating them onto specific regions indicated by $Mask_{3D}$. For each region in $Mask_{3D}$, we simply count the annotations and distribute their mode within the region.

Expansion of Reliable Pseudo Labels. The reliable pseudo labels \tilde{Y}^r are also propagated to $Mask_{3D}$ by selecting the mode of them in each region specified by each mask. If the proportion of \tilde{Y}^r that equals the mode $label_m$ exceeds a certain threshold η , they are considered as the labels for all points in the region. Here, η is a hyperparameter that determines whether the mode of reliable pseudo labels is expanded to the entire mask.

Finally, the expanded reliable pseudo labels are fused with the expanded annotations, denoted as \tilde{Y} . This procedure is shown in Algorithm 1.

Algorithm 1 Label Propagation

```

PROPAGATE( $Y, \tilde{Y}^r, Mask_{3D}, T, \eta$ )
  for  $t = 1$  to  $T$  do
     $mask \leftarrow Mask_{3D}(t)$ 
     $label_m \leftarrow mode(\tilde{Y}^r \cdot mask)$ 
    if  $\frac{|\{(\tilde{Y}^r \cdot mask) = label_m\}|}{|mask|} > \eta$  then
       $(\tilde{Y} \cdot mask) \leftarrow label_m$ 
    else if not_empty( $Y \cdot mask$ ) then
       $(\tilde{Y} \cdot mask) \leftarrow mode(Y \cdot mask)$ 
  return  $\tilde{Y}$ 

```

Even with the expansion of the available labels, they remain susceptible to inaccuracies or incompleteness because of misalignments among categories in the 2D segmentation mask,

inaccuracies in the 2D-3D projection, or points falling on boundaries and being erroneously propagated to nearby masks. This problem can be effectively tackled by implementing robust learning methods specifically designed to handle noisy labels.

D. Overall Loss Function

The propagated labels \tilde{Y} are treated as noisy labels, and the normalized loss [24] is applied to them, which are Normalized Cross-Entropy (NCE) loss and Reverse Cross-Entropy (RCE) loss:

$$\mathbf{L}_m = NCE[\tilde{Y}, P] + RCE[\tilde{Y}, P], \quad (10)$$

where P is the prediction of original point clouds X , and NCE and RCE are defined as:

$$\begin{aligned} NCE &= \frac{-\sum_{k=1}^K \mathbf{q}(k|x) \log \mathbf{p}(k|x)}{-\sum_{j=1}^K \sum_{k=1}^K \mathbf{q}(y=j|x) \log \mathbf{p}(k|x)}, \\ RCE &= -\sum_{k=1}^K \mathbf{p}(k|x) \log \mathbf{q}(k|x), \end{aligned} \quad (11)$$

where y is the label of input x in the K -class classification problem, $\mathbf{p}(k|x)$ is the softmax prediction, $\mathbf{q}(k|x)$ represents the distribution over different labels for sampler.

Furthermore, the Cross-Entropy (CE) loss is applied to the prediction P and the original sparse labels Y .

$$\mathbf{L}_{seg} = CE[Y, P] \quad (12)$$

Finally, the overall loss function is the weighted sum of the above losses:

$$\mathbf{L} = \lambda_{seg} \mathbf{L}_{seg} + \lambda_r \mathbf{L}_r + \lambda_a \mathbf{L}_a + \lambda_m \mathbf{L}_m \quad (13)$$

IV. EXPERIMENT

A. Experimental Setup

The experiments were conducted using two datasets: ScanNetV2 [48] and S3DIS [49]. ScanNetV2 comprises 1201 training and 312 validation scans sourced from 706 distinct scenes, with an additional test set of 100 scans. The experiments focus on the ‘3D Semantic Label with Limited Annotations’ benchmark, employing 20 training points per scene.

The S3DIS dataset includes 3D scans of 271 rooms across 6 distinct areas, each containing 13 categories. For training, data from Areas 1, 2, 3, 4, and 6 are utilized, and the model’s performance is evaluated on Area 5. Since there’s no official limited annotation setting provided, the ‘One-Thing-One-Click (OTOC)’ annotation scheme [13] is adopted by randomly retaining one annotation in each instance.

The evaluation metrics employed in this study involve the calculation of the mean of class-wise intersection over union (mIoU).

B. Comparison on the ScanNetV2 and S3DIS dataset

Table I shows the mIoU results of previous methods and our method on the ScanNetV2 3D semantic label benchmark. Table II presents the mIoU results obtained by previous methods alongside the proposed approach on the S3DIS 3D semantic segmentation dataset, with evaluations conducted on Area 5 as the testing set. For the two tables, the first part includes pure 3D data-based fully supervised methods, the second part shows fully supervised methods that leverage 2D information, and the third part includes weakly supervised methods based on 3D data. Our experiments are conducted on Point Transformer V3 (PTv3) [25] backbone, and results are shown in the fourth part.

On both datasets, our approach surpasses the baseline by more than 9%, outperforms some fully supervised methods, and achieves state-of-the-art performance among weakly supervised methods, demonstrating its effectiveness.

TABLE I
COMPARISON ON SCANNetV2 TESTING SET. OUR METHOD ACHIEVES THE TOP PERFORMANCE AMONG WEAKLY SUPERVISED METHODS, WITH 20 TRAINING POINTS PER SCENE.

Method	Supervision	mIoU(%)
PointNet++ [2]	100%	33.9
PointCNN [3]	100%	45.8
SparseConvNet [8]	100%	72.5
KPConv [4]	100%	68.6
MinkowskiNet [9]	100%	73.6
PointConv [5]	100%	66.6
Point Transformer V2 [7]	100%	75.2
3DMV [50]	100%	48.4
Virtual MVFusion [19]	100%	74.6
BPNet [21]	100%	74.9
PointContrast_LA_SEM [51]	20 points	55.0
Viewpoint_BN_LA_AIR [52]	20 points	54.8
One-Thing-One-Click [13]	20 points	59.4
PointMatch [17]	20 points	62.4
RAC-Net [18]	20 points	63.9
Our Baseline (PTv3)	20 points	60.1
Ours	20 points	69.4
Our Upper Bound (PTv3)	100%	77.9

TABLE II
COMPARISON ON S3DIS TESTING SET (AREA 5). OUR METHOD ACHIEVES THE BEST RESULTS AMONG WEAKLY SUPERVISED METHODS UNDER THE ‘OTOC’ SETTING.

Method	Supervision	mIoU(%)
PointNet [1]	100%	41.1
PointCNN [3]	100%	57.3
KPConv [4]	100%	65.4
MinkowskiNet [9]	100%	65.4
Point Transformer V2 [7]	100%	71.6
Virtual MVFusion [19]	100%	65.4
One-Thing-One-Click [13]	0.02%(OTOC)	50.1
PointMatch [17]	0.02%(OTOC)	55.3
RAC-Net [18]	0.02%(OTOC)	58.4
Our Baseline (PTv3)	0.02%(OTOC)	54.7
Ours	0.02%(OTOC)	64.4
Our Upper Bound (PTv3)	100%	73.4

Our method can be applied with any point cloud segmentation backbone. We conducted experiments using the ScanNetV2 dataset with the SparseUNet [53] backbone to validate its effectiveness. Table III presents the results of the ScanNet V2 validation set.

TABLE III
COMPARISON ON SCANNetV2 VALIDATION SET WITH SPARSEUNET [53] BACKBONE.

Method	Supervision	mIoU(%)
Our Baseline (SparseUNet)	20 points	59.2
Ours	20 points	67.4
Our Upper Bound (SparseUNet)	100%	74.8

C. Discussions

1) *Effectiveness of Back-Projected Masks*: This section explores the effectiveness of back-projecting 2D segmentation masks to 3D masks and propagating the limited annotations Y onto the masks, by comparing the number and accuracy of expanded labels and training the model using the expanded labels.

Number and Accuracy of Expanded Annotations. Table IV presents the average number and accuracy of labels within a scene on the ScanNetV2 dataset and the S3DIS dataset, both before and after the process of propagating the initial limited annotations onto the back-projected 3D masks. This demonstrates the efficacy of leveraging 2D images to significantly augment label coverage while maintaining a relatively high level of accuracy.

Visualization of Expanded Labels after Label Initialization. Some examples of expanded labels obtained by label

TABLE IV
THE AVERAGE NUMBER AND ACCURACY OF LABELS FROM SCANNetV2 AND S3DIS TRAINING SET, BEFORE AND AFTER LABEL INITIALIZATION. FOR THE SCANNetV2 DATASET, THE BASELINE IS 20 ANNOTATIONS PER SCENE, BUT THE AVERAGE IS 15 DUE TO SOME RANDOMLY RETAINED POINTS INITIALLY LACKING GROUND TRUTH. WE ALSO COMPARE THE PROPAGATION OF LIMITED ANNOTATIONS TO 3D MASKS GENERATED BY SAM3D [39]. FOR THE S3DIS DATASET, WE ADOPT THE ‘OTOC’ SETTING, RANDOMLY RETAINING 1 ANNOTATION PER INSTANCE, RESULTING IN AN AVERAGE OF 36.7 ANNOTATIONS PER SCENE.

Dataset	Scheme	Number	Accuracy
ScanNetV2	Fully Supervised Annotations	145170.8	100%
	Limited Annotations	15.0	100%
	Expanded Labels (Ours)	5565.5	93.0%
	Expanded Labels (SAM3D)	6224.5	81.8%
S3DIS	Fully Supervised Annotations	955036.4	100%
	‘OTOC’ Annotations	36.7	100%
	Expanded Labels (Ours)	115732.7	71.6%

initialization stage are displayed in Figure 2 and Figure 3 for ScanNetV2 and S3DIS training set, showing the original point clouds, 3D segmentation masks, initial sparse annotations, expanded labels onto the masks after label initialization, and the ground truth. The expanded labels are mostly accurate compared to the ground truth, but they still cluster around the initial sparse labels. Therefore, in subsequent steps, we select reliable pseudo labels and propagate them onto the masks to fully utilize these masks.

Figure 4 illustrates a failure case in label initialization, where a point falling on the boundary is chosen as one of the limited annotations. This label is propagated to another masked region, filling it and introducing noise into the expanded labels.

Utilizing 3D Masks from SAM3D. We also explored using segmentation masks generated by SAM3D [39] for 2D images

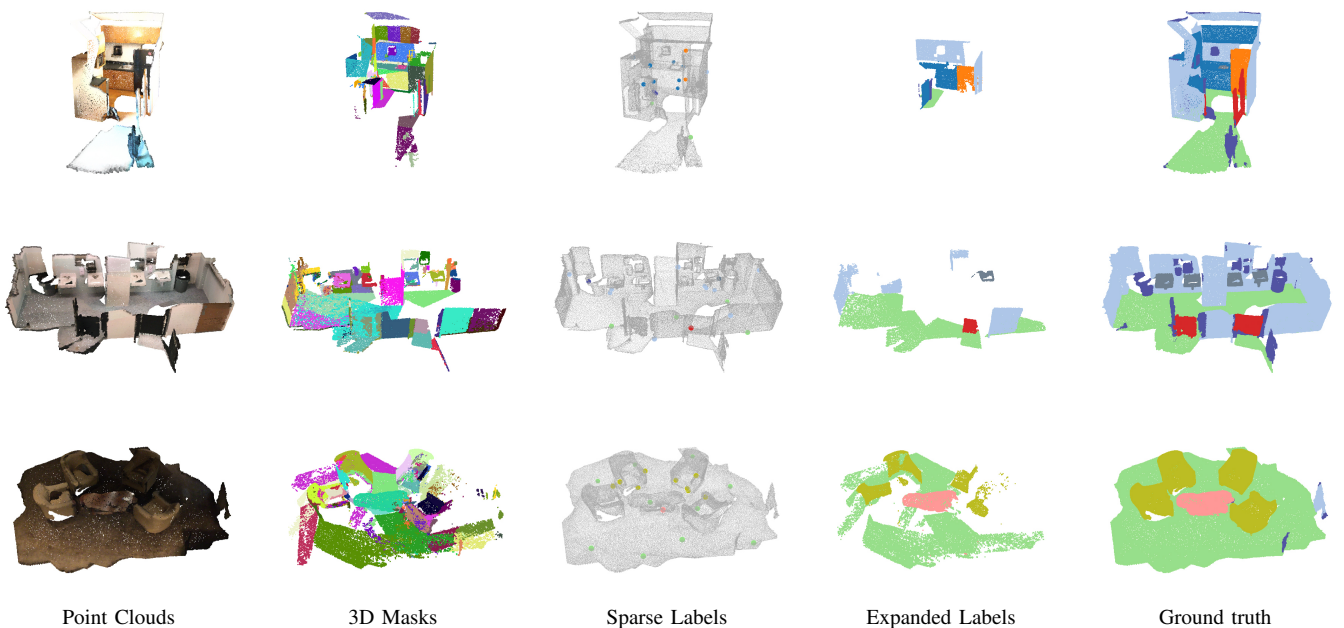


Fig. 2. Visualization of 3D labels from ScanNetV2 dataset before and after label initialization, with 20 labeled points per scene. The 3D masks projected from the 2D masks accurately represent the position and shape of the 3D objects. The expanded labels are generally accurate compared to the ground truth but do not cover all the masked regions.

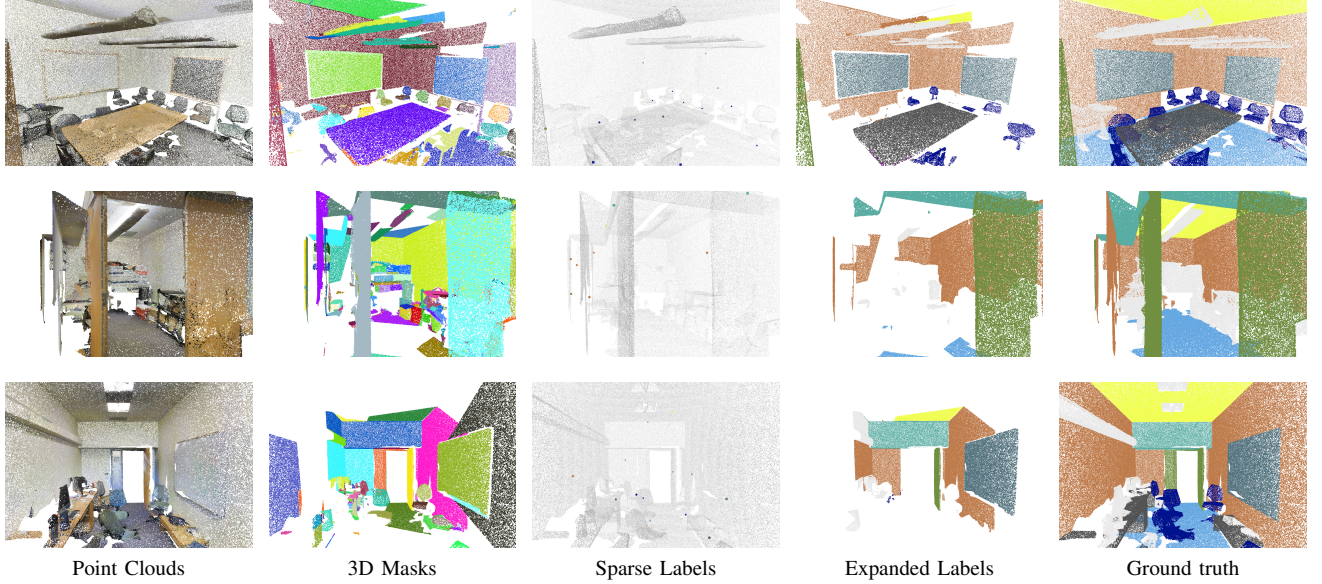


Fig. 3. Visualization of 3D labels from S3DIS dataset before and after label initialization with ‘OTOC’ annotation scheme.

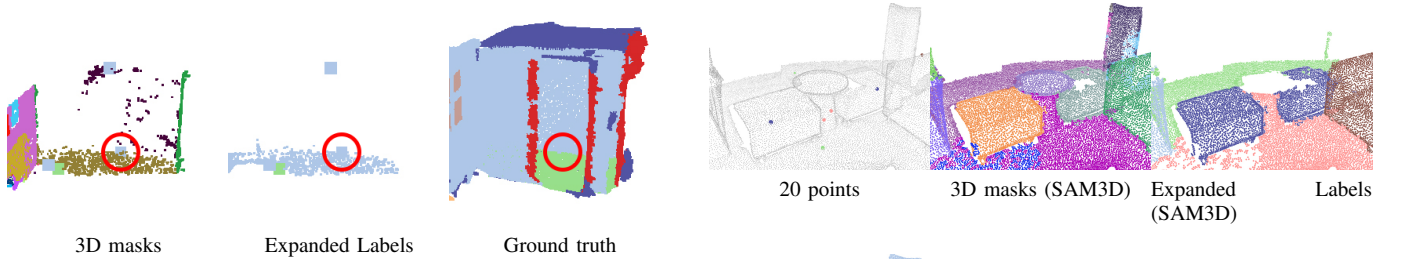


Fig. 4. A failure case from the ScanNetV2 dataset after label initialization. The circled point is one of the initial limited annotations, labeled as the wall ■. However, being located on the boundary between the floor and wall masks, it incorrectly propagates to the floor ■.

corresponding to each 3D point cloud, as a replacement for our generated masks, and incorporated them into our model’s training. SAM3D, by leveraging all RGB images in each scene, produces more complete masks but is more time-consuming. However, it sometimes merges masks of different objects into a single mask, introducing additional noise when limited annotations are propagated onto the masks, as shown in Table IV and Figure 5. Using masks generated by SAM3D, we trained the model on the ScanNetV2 training set with Point Transformer V3 as the backbone and compared it with our methods on the validation set in Table V.

TABLE V
COMPARISON ON SCANNETV2 VALIDATION SET BY TRAINING MODELS
USING 3D MASKS GENERATED BY OUR METHODS AND SAM3D.

Method	Supervision	mIoU(%)
Our Baseline (PTv3)	20 points	60.6
SAM3D	20 points	70.2
Ours	20 points	71.3
Our Upper Bound (PTv3)	100%	77.5

2) *Noise-Robust Loss on the Expanded Labels*: The expanded sparse annotations can be utilized directly by in-

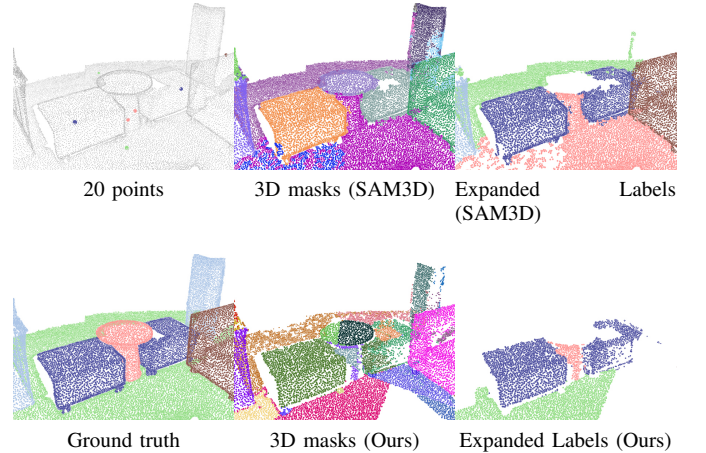


Fig. 5. A failure case occurred when using SAM3D for label initialization, where two masks were incorrectly merged, causing both the floor ■ and table ■ to be treated as table ■ when the initial sparse annotations were propagated onto the mask.

corporating an additional loss term alongside the Cross-Entropy loss for original sparse labels utilized in the baseline method, without the need for additional techniques, i.e. $\mathbf{L} = \lambda_{\text{seg}}\mathbf{L}_{\text{seg}} + \lambda_{\text{m}}\mathbf{L}_{\text{m}}$.

Because reality lacks ground truth, and one point may not adequately represent the entire mask due to noise in projected masks or misalignment among classes, the accuracy of expanded labels is uncertain. Therefore, they are regarded as noisy labels. As a result, noise-robust loss functions [24], [44] are compared on the ScanNetV2 dataset on Point Transformer V3 backbone, as shown in Table VI. Moreover, Table VII shows the results on the S3DIS dataset. These results demonstrate that despite the noise in the expanded labels, their direct application with an additional loss can enhance performance. This also proves the efficacy of propagating limited annotations onto back-projected masks.

TABLE VI

RESULTS ON SCANNetV2 VALIDATION SET WITH THE NOISE-ROBUST LOSS ON EXPANDED SPARSE ANNOTATIONS. IN PARTICULAR, THE ‘NCE’ AND ‘RCE’ LOSSES ACHIEVED THE HIGHEST RESULTS, SO WE CHOSE TO APPLY THEM TO ALL LABELS EXPANDED TO 3D MASKS.

CE	NFL	MAE	RCE	NCE	AGCE	AUE	mIoU (%)
-	-	-	-	-	-	-	60.6
✓	-	-	-	-	-	-	62.4
-	✓	✓	-	-	-	-	60.8
-	✓	-	✓	-	-	-	62.8
-	-	✓	-	✓	-	-	60.9
-	-	-	✓	✓	-	-	63.0
-	-	-	-	✓	✓	-	62.3
-	-	-	-	✓	-	✓	61.4

TABLE VII

RESULTS ON S3DIS TESTING SET WITH NOISE-ROBUST LOSS ON EXPANDED SPARSE ANNOTATIONS.

CE	NFL	MAE	RCE	NCE	AGCE	AUE	mIoU (%)
-	-	-	-	-	-	-	54.7
✓	-	-	-	-	-	-	56.3
-	-	-	✓	✓	-	-	57.1

3) *Selection of η* : As mentioned in Section III-C, for reliable pseudo labels, due to their potential inaccuracy or instability, it is desirable for them to occupy a certain proportion of each mask-represented region, and then expand them to the entire mask. The hyperparameter η is introduced to control this proportion.

Comparison on different values of η . On the one hand, it is desirable to expand reliable pseudo labels as much as possible; on the other hand, the higher the proportion they occupy on the mask, the more reliable the expanded labels are. Table VIII demonstrates the effects of training the model with different values of η .

TABLE VIII

COMPARISON ON DIFFERENT VALUE OF η ON SCANNetV2 VALIDATION SET.

η_0	mIoU (%)
0.3	70.2
0.5	70.4
0.7	71.3
0.9	70.7

Number and Accuracy of Expanded Labels \tilde{Y} . Figure 6 shows the number and accuracy of the expanded limited annotations and reliable pseudo labels after propagating them onto the masks, defined as \tilde{Y} . From the figure, it can be observed that the accuracy of \tilde{Y} tends to stabilize in the later stages, with its quantity increasing with more epochs.

4) *Ablations*: In this section, we assess the effectiveness of each loss function and the corresponding computational entities within the total objective function on the ScanNetV2 dataset, with 20 training points in each scene. The results of the ScanNetV2 validation set are shown in Table IX. These results demonstrate the effectiveness of the different losses and their components.

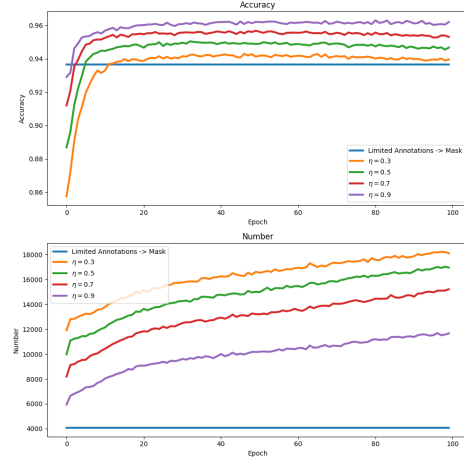


Fig. 6. Number and accuracy of \tilde{Y} with different value of η through the training process.

Corresponding to these computational entities, Figure 7 shows an example of the expanded labels and ground truth. This demonstrates that propagating the initial sparse annotations to the masks and subsequently adding reliable pseudo labels to the masks can expand the available labels with high accuracy.

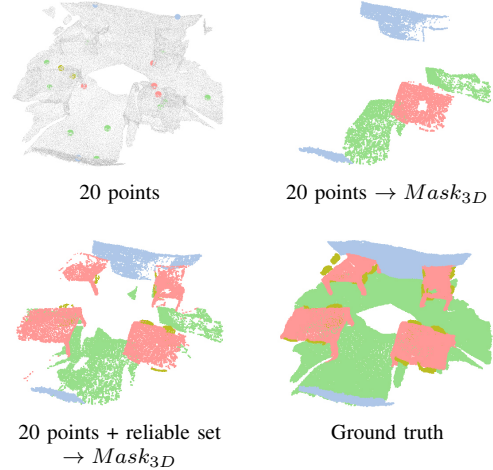


Fig. 7. An example of expanded labels from the ScanNetV2 dataset. The available labels are greatly expanded after propagating initial annotations and reliable pseudo labels to the masks.

5) *Result Visualization*: This section qualitatively presents the results of our model on ScanNetV2 and S3DIS datasets, as shown in Figures 8 and 9. Compared to ground truth, the model is capable of accurately generating semantic segmentation masks. However, there is a possibility of erroneously grouping small objects with nearby ones, or encountering some ambiguity at object boundaries.

V. CONCLUSION

In this paper, we introduce an innovative approach for 3D weakly supervised semantic segmentation by integrating 2D images. Our method aims to effectively augment the sparse

TABLE IX
ABLATIONS ON SCANNETV2 VALIDATION SET.

CE Loss	NCEandRCE	CE Loss	KL Loss	mIoU(%)
20 points	-	-	-	60.6
20 points	20 points \rightarrow $Mask_{3D}$	-	-	63.0
20 points	-	reliable set	ambiguous set	69.9
20 points	20 points \rightarrow $Mask_{3D}$	reliable set	ambiguous set	70.8
20 points	reliable set \rightarrow $Mask_{3D}$	reliable set	ambiguous set	70.4
20 points	20 points + reliable set \rightarrow $Mask_{3D}$	reliable set	ambiguous set	71.3

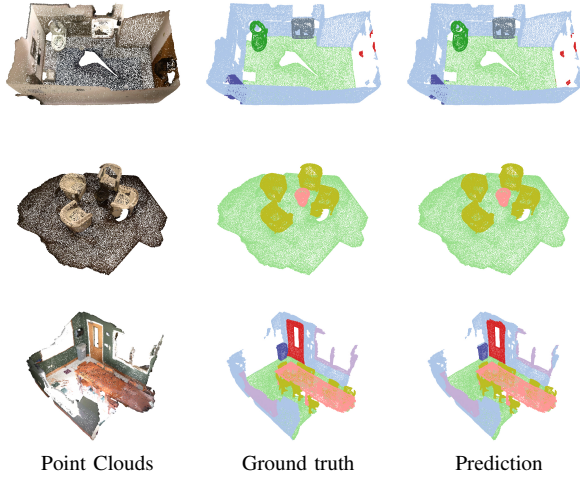


Fig. 8. Examples of semantic segmentation results on ScanNetV2 validation set. Our model can perform semantic segmentation with high accuracy, but there are still some issues with edges and small objects.

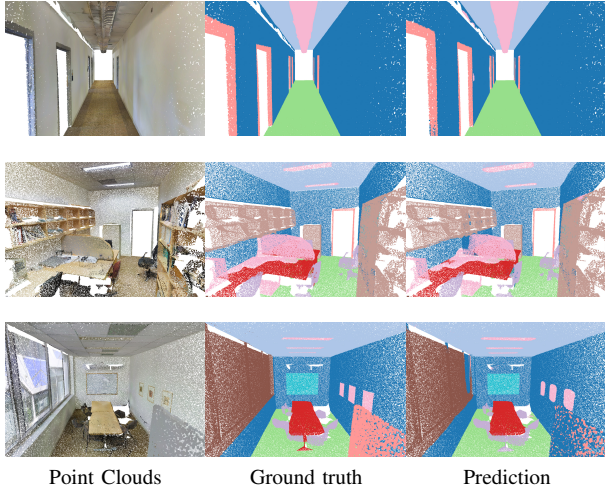


Fig. 9. Examples of semantic segmentation results on S3DIS testing set (Area 5).

labeling of 3D point clouds by leveraging the geometric correspondence between 2D views and 3D point clouds. This involves utilizing segmentation masks derived from 2D foundational models and back-projecting them to 3D space. By propagating the initial limited annotations onto the 3D masks, we substantially increase the available labels. Additionally, we incorporate consistency regularization and select

reliable pseudo labels, which are then proportionally expanded into 3D masks, maximizing their utility. To address noise in the expanded labels, we employ noise-robust techniques to enhance model performance. Experiments on ScanNetV2 and S3DIS datasets demonstrate state-of-the-art performance.

For future works, there is still room for improvement in the model's performance because issues like unclear boundaries and insensitivity to small objects persist. Besides, broader applications of 3D masks can be explored to maximize their effectiveness. Exploring alternative methods for handling label noise is also necessary. In the future, our method could be further explored for applications in various fields, including outdoor datasets, such as those collected by autonomous driving systems, where images and point clouds are acquired simultaneously.

REFERENCES

- [1] R. Q. Charles, H. Su, M. Kaichun, and L. J. Guibas, "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Honolulu, HI: IEEE, Jul. 2017, pp. 77–85. [Online]. Available: <http://ieeexplore.ieee.org/document/8099499/>
- [2] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space," Jun. 2017, arXiv:1706.02413 [cs]. [Online]. Available: <http://arxiv.org/abs/1706.02413>
- [3] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, "PointCNN: Convolution on X-Transformed Points," in *Advances in Neural Information Processing Systems*, vol. 31. Curran Associates, Inc., 2018. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2018/hash/f5f8590cd58a54e94377e6ae2eded4d9-Abstract.html
- [4] H. Thomas, C. R. Qi, J.-E. Deschaud, B. Marcotegeui, F. Goulette, and L. J. Guibas, "KPConv: Flexible and Deformable Convolution for Point Clouds," Aug. 2019, arXiv:1904.08889 [cs]. [Online]. Available: <http://arxiv.org/abs/1904.08889>
- [5] W. Wu, Z. Qi, and L. Fuxin, "PointConv: Deep Convolutional Networks on 3D Point Clouds," Nov. 2020, arXiv:1811.07246 [cs]. [Online]. Available: <http://arxiv.org/abs/1811.07246>
- [6] H. Zhao, L. Jiang, J. Jia, P. Torr, and V. Koltun, "Point Transformer," Sep. 2021, arXiv:2012.09164 [cs]. [Online]. Available: <http://arxiv.org/abs/2012.09164>
- [7] X. Wu, Y. Lao, L. Jiang, X. Liu, and H. Zhao, "Point Transformer V2: Grouped Vector Attention and Partition-based Pooling," Oct. 2022, arXiv:2210.05666 [cs]. [Online]. Available: <http://arxiv.org/abs/2210.05666>
- [8] B. Graham, M. Engelcke, and L. V. D. Maaten, "3D Semantic Segmentation with Submanifold Sparse Convolutional Networks," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Salt Lake City, UT, USA: IEEE, Jun. 2018, pp. 9224–9232. [Online]. Available: <https://ieeexplore.ieee.org/document/8579059/>
- [9] C. Choy, J. Gwak, and S. Savarese, "4D Spatio-Temporal ConvNets: Minkowski Convolutional Neural Networks," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Long Beach, CA, USA: IEEE, Jun. 2019, pp. 3070–3079. [Online]. Available: <https://ieeexplore.ieee.org/document/8953494/>

- [10] B. Wu, A. Wan, X. Yue, and K. Keutzer, "SqueezeSeg: Convolutional Neural Nets with Recurrent CRF for Real-Time Road-Object Segmentation from 3D LiDAR Point Cloud," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. Brisbane, QLD: IEEE, May 2018, pp. 1887–1893. [Online]. Available: <https://ieeexplore.ieee.org/document/8462926/>
- [11] Z. J. Yew and G. H. Lee, "3DFeat-Net: Weakly Supervised Local 3D Features for Point Cloud Registration," in *Computer Vision – ECCV 2018*, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds. Cham: Springer International Publishing, 2018, vol. 11219, pp. 630–646, series Title: Lecture Notes in Computer Science. [Online]. Available: https://link.springer.com/10.1007/978-3-030-01267-0_37
- [12] J. Wei, G. Lin, K.-H. Yap, T.-Y. Hung, and L. Xie, "Multi-Path Region Mining for Weakly Supervised 3D Semantic Segmentation on Point Clouds," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Seattle, WA, USA: IEEE, Jun. 2020, pp. 4383–4392. [Online]. Available: <https://ieeexplore.ieee.org/document/9157503/>
- [13] Z. Liu, X. Qi, and C.-W. Fu, "One Thing One Click: A Self-Training Approach for Weakly Supervised 3D Semantic Segmentation," in *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Nashville, TN, USA: IEEE, Jun. 2021, pp. 1726–1736. [Online]. Available: <https://ieeexplore.ieee.org/document/9578763/>
- [14] Z. Wu, Y. Wu, G. Lin, J. Cai, and C. Qian, "Dual Adaptive Transformations for Weakly Supervised Point Cloud Segmentation," Jul. 2022, arXiv:2207.09084 [cs]. [Online]. Available: <http://arxiv.org/abs/2207.09084>
- [15] A. Tao, Y. Duan, Y. Wei, J. Lu, and J. Zhou, "SegGroup: Seg-Level Supervision for 3D Instance and Semantic Segmentation," *IEEE Transactions on Image Processing*, vol. 31, pp. 4952–4965, 2022, arXiv:2012.10217 [cs]. [Online]. Available: <http://arxiv.org/abs/2012.10217>
- [16] Z. Liu, X. Qi, and C.-W. Fu, "One Thing One Click++: Self-Training for Weakly Supervised 3D Scene Understanding," Mar. 2023, arXiv:2303.14727 [cs]. [Online]. Available: <http://arxiv.org/abs/2303.14727>
- [17] Y. Wu, Z. Yan, S. Cai, G. Li, X. Han, and S. Cui, "PointMatch: A consistency training framework for weakly supervised semantic segmentation of 3D point clouds," *Computers & Graphics*, vol. 116, pp. 427–436, Nov. 2023. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S0097849323002297>
- [18] Z. Wu, Y. Wu, G. Lin, and J. Cai, "Reliability-Adaptive Consistency Regularization for Weakly-Supervised Point Cloud Segmentation," *International Journal of Computer Vision*, Jan. 2024. [Online]. Available: <https://link.springer.com/10.1007/s11263-023-01975-8>
- [19] A. Kundu, X. Yin, A. Fathi, D. Ross, B. Brewington, T. Funkhouser, and C. Pantofaru, "Virtual Multi-view Fusion for 3D Semantic Segmentation," Jul. 2020, arXiv:2007.13138 [cs, eess]. [Online]. Available: <http://arxiv.org/abs/2007.13138>
- [20] I. Alonso, L. Riazuelo, L. Montesano, and A. C. Murillo, "3D-MiniNet: Learning a 2D Representation From Point Clouds for Fast and Efficient 3D LIDAR Semantic Segmentation," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5432–5439, Oct. 2020. [Online]. Available: <https://ieeexplore.ieee.org/document/9134888/>
- [21] W. Hu, H. Zhao, L. Jiang, J. Jia, and T.-T. Wong, "Bidirectional Projection Network for Cross Dimension Scene Understanding," Mar. 2021, arXiv:2103.14326 [cs]. [Online]. Available: <http://arxiv.org/abs/2103.14326>
- [22] K. Genova, X. Yin, A. Kundu, C. Pantofaru, F. Cole, A. Sud, B. Brewington, B. Shucker, and T. Funkhouser, "Learning 3D Semantic Segmentation with only 2D Image Supervision," in *2021 International Conference on 3D Vision (3DV)*. London, United Kingdom: IEEE, Dec. 2021, pp. 361–372. [Online]. Available: <https://ieeexplore.ieee.org/document/9665849/>
- [23] F. Li, H. Zhang, P. Sun, X. Zou, S. Liu, J. Yang, C. Li, L. Zhang, and J. Gao, "Semantic-SAM: Segment and Recognize Anything at Any Granularity," Jul. 2023, arXiv:2307.04767 [cs]. [Online]. Available: <http://arxiv.org/abs/2307.04767>
- [24] X. Ma, H. Huang, Y. Wang, S. Romano, S. Erfani, and J. Bailey, "Normalized Loss Functions for Deep Learning with Noisy Labels," Jun. 2020, arXiv:2006.13554 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/2006.13554>
- [25] X. Wu, L. Jiang, P.-S. Wang, Z. Liu, X. Liu, Y. Qiao, W. Ouyang, T. He, and H. Zhao, "Point Transformer V3: Simpler, Faster, Stronger," Mar. 2024, arXiv:2312.10035 [cs]. [Online]. Available: <http://arxiv.org/abs/2312.10035>
- [26] G. Liu, O. Van Kaick, H. Huang, and R. Hu, "Active self-training for weakly supervised 3D scene semantic segmentation," *Computational Visual Media*, vol. 10, no. 3, pp. 425–438, Jun. 2024. [Online]. Available: <https://link.springer.com/10.1007/s41095-022-0311-7>
- [27] S. Dong, F. Liu, and G. Lin, "Leveraging Large-Scale Pretrained Vision Foundation Models for Label-Efficient 3D Point Cloud Segmentation," Nov. 2023, arXiv:2311.01989 [cs]. [Online]. Available: <http://arxiv.org/abs/2311.01989>
- [28] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick, "Segment Anything," Apr. 2023, arXiv:2304.02643 [cs]. [Online]. Available: <http://arxiv.org/abs/2304.02643>
- [29] W. Liu, C. Zhang, G. Lin, and F. Liu, "Crnet: Cross-reference networks for few-shot segmentation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 4165–4173.
- [30] W. Liu, C. Zhang, H. Ding, T.-Y. Hung, and G. Lin, "Few-shot segmentation with optimal transport matching and message flow," *IEEE Transactions on Multimedia*, vol. 25, pp. 5130–5141, 2022.
- [31] W. Liu, Z. Wu, Y. Zhao, Y. Fang, C.-S. Foo, J. Cheng, and G. Lin, "Harmonizing base and novel classes: A class-contrastive approach for generalized few-shot segmentation," *International Journal of Computer Vision*, vol. 132, no. 4, pp. 1277–1291, 2024.
- [32] W. Liu, W. Zhou, J. Liu, P. Hu, J. Cheng, J. Han, and W. Lin, "Modality-aware feature matching: A comprehensive review of single-and cross-modality techniques," *arXiv preprint arXiv:2507.22791*, 2025.
- [33] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Boston, MA, USA: IEEE, Jun. 2015, pp. 3431–3440. [Online]. Available: <http://ieeexplore.ieee.org/document/7298965/>
- [34] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," May 2015, arXiv:1505.04597 [cs]. [Online]. Available: <http://arxiv.org/abs/1505.04597>
- [35] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. S. Torr, "Conditional Random Fields as Recurrent Neural Networks," in *2015 IEEE International Conference on Computer Vision (ICCV)*. Santiago, Chile: IEEE, Dec. 2015, pp. 1529–1537. [Online]. Available: <http://ieeexplore.ieee.org/document/7410536/>
- [36] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," Jun. 2021, arXiv:2010.11929 [cs]. [Online]. Available: <http://arxiv.org/abs/2010.11929>
- [37] B. Zhang, Z. Tian, Q. Tang, X. Chu, X. Wei, C. Shen, and Y. Liu, "SegViT: Semantic Segmentation with Plain Vision Transformers," Dec. 2022, arXiv:2210.05844 [cs]. [Online]. Available: <http://arxiv.org/abs/2210.05844>
- [38] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. Montreal, QC, Canada: IEEE, Oct. 2021, pp. 9992–10002. [Online]. Available: <https://ieeexplore.ieee.org/document/9710580/>
- [39] Y. Yang, X. Wu, T. He, H. Zhao, and X. Liu, "SAM3D: Segment Anything in 3D Scenes," Jun. 2023, arXiv:2306.03908. [Online]. Available: <http://arxiv.org/abs/2306.03908>
- [40] X. Wang, S. Li, K. Kallidromitis, Y. Kato, K. Kozuka, and T. Darrell, "Hierarchical Open-vocabulary Universal Image Segmentation," Jul. 2023, arXiv:2307.00764 [cs]. [Online]. Available: <http://arxiv.org/abs/2307.00764>
- [41] A. Ghosh, H. Kumar, and P. S. Sastry, "Robust Loss Functions under Label Noise for Deep Neural Networks," Dec. 2017, arXiv:1712.09482 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1712.09482>
- [42] Z. Zhang and M. R. Sabuncu, "Generalized Cross Entropy Loss for Training Deep Neural Networks with Noisy Labels," Nov. 2018, arXiv:1805.07836 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1805.07836>
- [43] Y. Wang, X. Ma, Z. Chen, Y. Luo, J. Yi, and J. Bailey, "Symmetric Cross Entropy for Robust Learning with Noisy Labels," Aug. 2019, arXiv:1908.06112 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1908.06112>
- [44] X. Zhou, X. Liu, J. Jiang, X. Gao, and X. Ji, "Asymmetric Loss Functions for Learning with Noisy Labels," Jun. 2021, arXiv:2106.03110 [cs]. [Online]. Available: <http://arxiv.org/abs/2106.03110>

- [45] E. Arazo, D. Ortego, P. Albert, N. E. O'Connor, and K. McGuinness, "Unsupervised Label Noise Modeling and Loss Correction," Jun. 2019, arXiv:1904.11238 [cs]. [Online]. Available: <http://arxiv.org/abs/1904.11238>
- [46] S. Ye, D. Chen, S. Han, and J. Liao, "Learning with Noisy Labels for Robust Point Cloud Segmentation," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. Montreal, QC, Canada: IEEE, Oct. 2021, pp. 6423–6432. [Online]. Available: <https://ieeexplore.ieee.org/document/9710169/>
- [47] S. Kim, S. Lee, D. Hwang, J. Lee, S. J. Hwang, and H. J. Kim, "Point Cloud Augmentation with Weighted Local Transformations," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. Montreal, QC, Canada: IEEE, Oct. 2021, pp. 528–537. [Online]. Available: <https://ieeexplore.ieee.org/document/9710410/>
- [48] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Niessner, "ScanNet: Richly-Annotated 3D Reconstructions of Indoor Scenes," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Honolulu, HI: IEEE, Jul. 2017, pp. 2432–2443. [Online]. Available: <https://ieeexplore.ieee.org/document/8099744/>
- [49] I. Armeni, S. Sax, A. R. Zamir, and S. Savarese, "Joint 2D-3D-Semantic Data for Indoor Scene Understanding," Apr. 2017, arXiv:1702.01105 [cs]. [Online]. Available: <http://arxiv.org/abs/1702.01105>
- [50] A. Dai and M. Nießner, "3DMV: Joint 3D-Multi-View Prediction for 3D Semantic Scene Segmentation," Mar. 2018, arXiv:1803.10409 [cs]. [Online]. Available: <http://arxiv.org/abs/1803.10409>
- [51] S. Xie, J. Gu, D. Guo, C. R. Qi, L. J. Guibas, and O. Litany, "PointContrast: Unsupervised Pre-training for 3D Point Cloud Understanding," Nov. 2020, arXiv:2007.10985 [cs]. [Online]. Available: <http://arxiv.org/abs/2007.10985>
- [52] L. Luo, B. Tian, H. Zhao, and G. Zhou, "Pointly-supervised 3D Scene Parsing with Viewpoint Bottleneck," Sep. 2021, arXiv:2109.08553 [cs]. [Online]. Available: <http://arxiv.org/abs/2109.08553>
- [53] S. Contributors, "Spconv: Spatially Sparse Convolution Library," 2022. [Online]. Available: <https://github.com/traveller59/spconv>